

# New Data Sources – Accessibility and Use

Frauke Kreuter

JPSM – Uni Mannheim – IAB

Paris 6/25/2019



## AAPOR Report on Big Data

AAPOR Big Data Task Force  
February 12, 2015

Prepared for AAPOR Council by the Task Force, with Task Force members including:

*Lilli Japac, Co-Chair, Statistics Sweden*  
*Franke Kreuter, Co-Chair, JPSM at the U. of Maryland, U. of Mannheim & IAB*  
*Marcus Berg, Stockholm University*  
*Paul Biemer, RTI International*  
*Paul Decker, Mathematica Policy Research*  
*Cliff Lampe, School of Information at the University of Michigan*  
*Julia Lane, American Institutes for Research*  
*Cathy O'Neil, Johnson Research Labs*  
*Abe Usher, HumanGeo Group*

Acknowledgement: We are grateful for comments, feedback and editorial help from Eran Ben-Porath, Jason McMillan, and the AAPOR council members.

The National Academies of  
SCIENCES • ENGINEERING • MEDICINE

REPORT

# INNOVATIONS IN FEDERAL STATISTICS

Combining Data Sources While  
Protecting Privacy

The National Academies of  
SCIENCES • ENGINEERING • MEDICINE

CONSENSUS STUDY REPORT

# FEDERAL STATISTICS, MULTIPLE DATA SOURCES, AND PRIVACY PROTECTION

Next Steps



SPONSORED BY THE



Federal Ministry  
of Education  
and Research

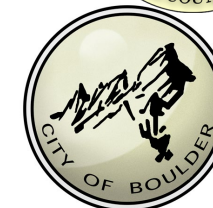
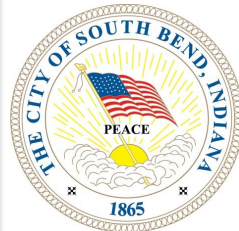
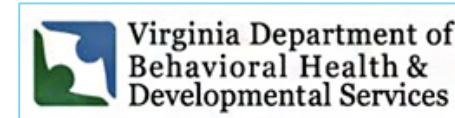
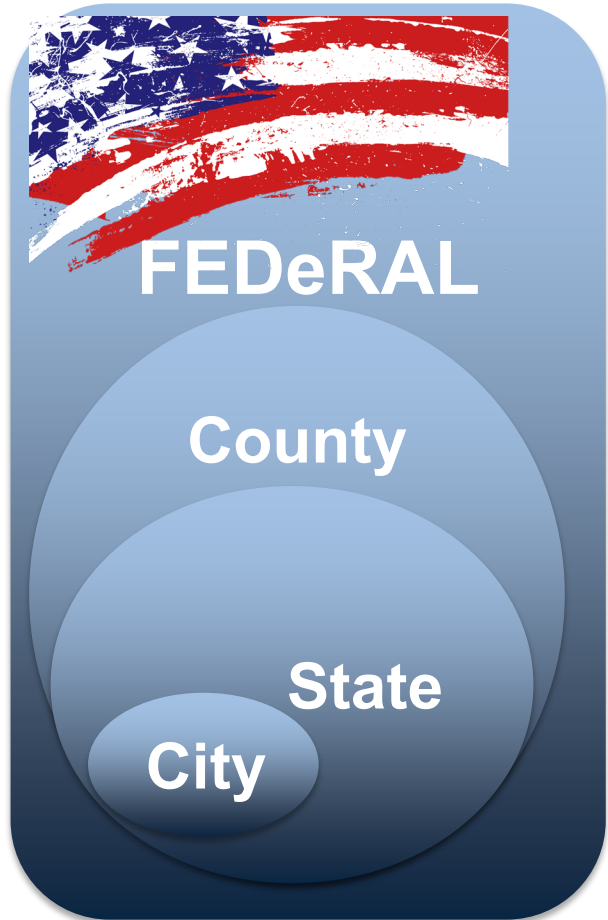


The project on which this report is based was funded by the Federal Ministry of Education and Research under the number [16OH22064]. Responsibility for the contents of this publication lies with the author.





# ~40 agencies from city, state, county and federal agencies



# The Hope

2

3 **Steuernummer**

4 eTIN lt. Lohnsteuerbescheinigung(en), sofern vorhanden  eTIN lt. weiterer Lohnsteuerbeschei

**Einkünfte aus nichtselbständiger Arbeit**

5 **Angaben zum Arbeitslohn** Lohnsteuerbescheinigung(en) Steuerklasse 1 - 5

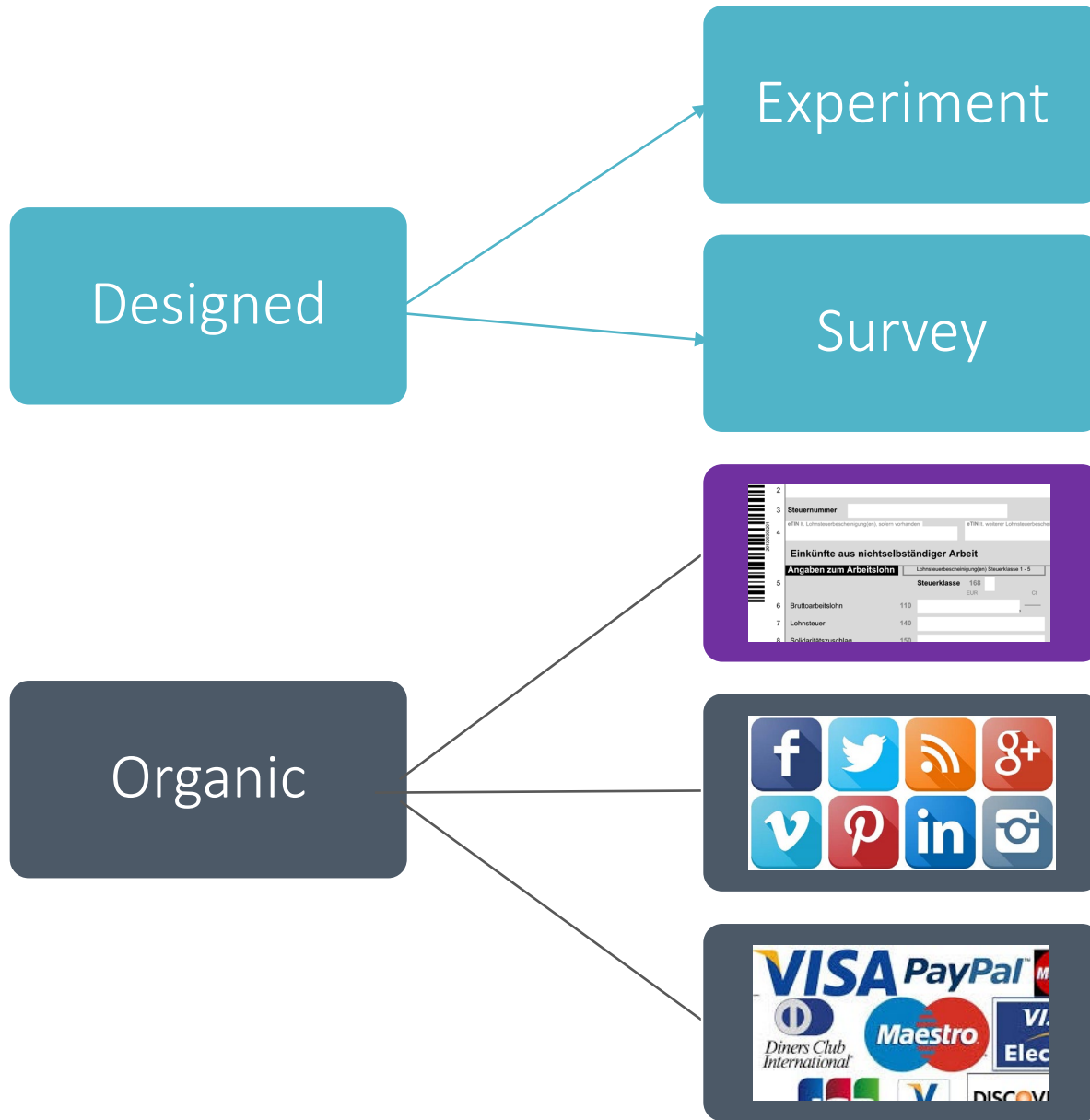
6 **Steuerklasse** 168  EUR  Ct

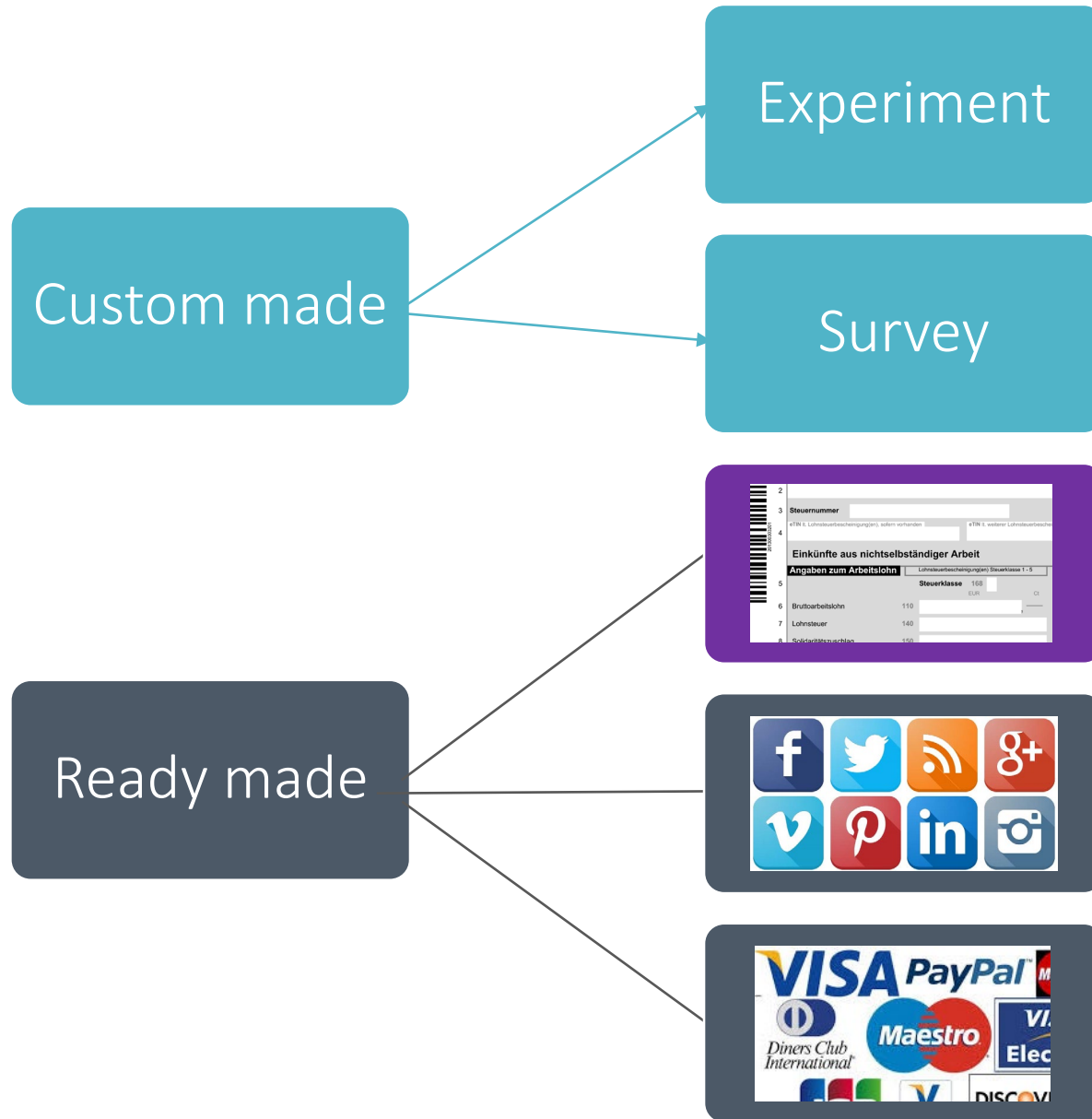
7 Bruttoarbeitslohn 110

8 Lohnsteuer 140

9 Solidaritätszuschlag 150







# Economic Indicators

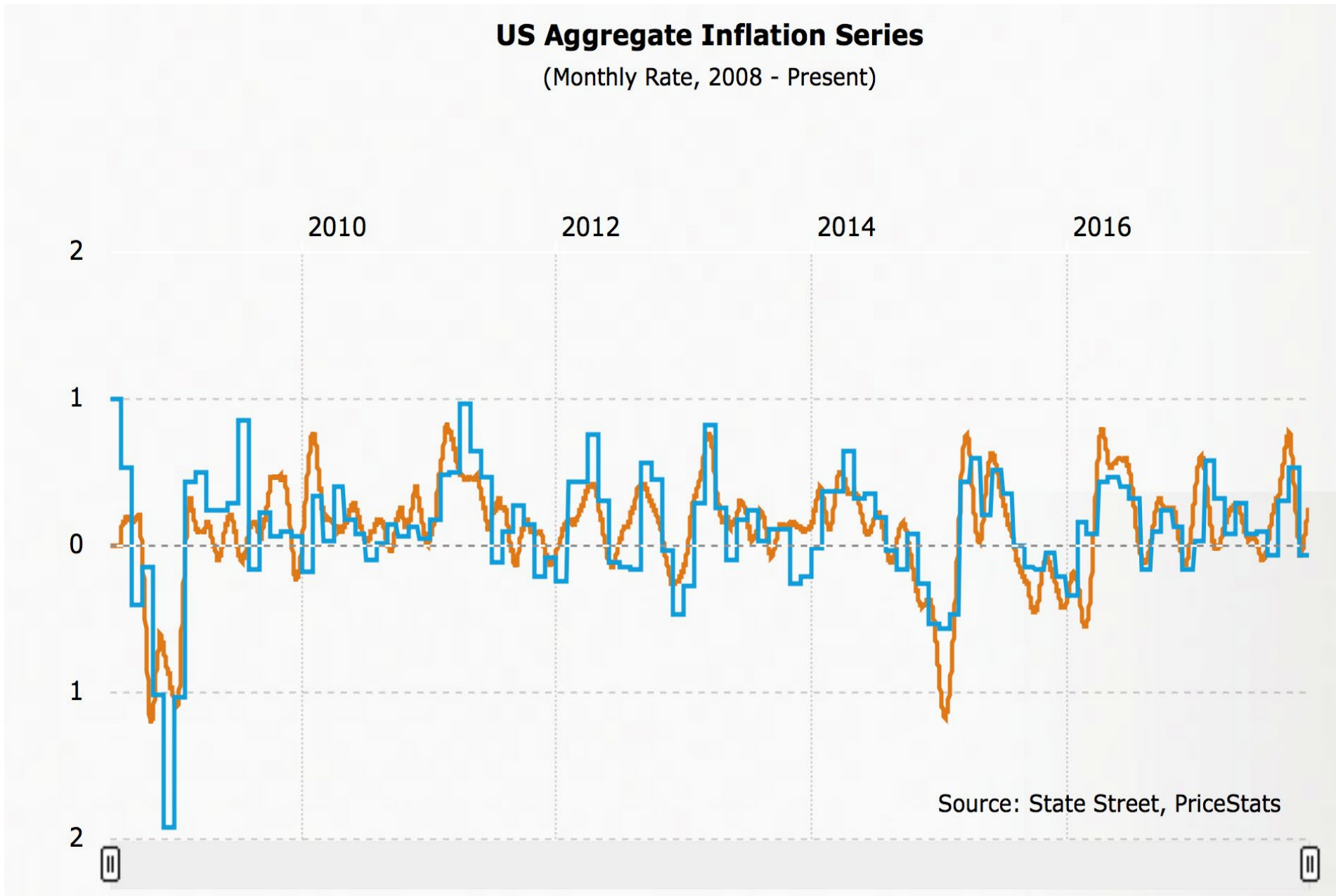
Examples for online data collection (and analysis)



# US Aggregate Inflation Series

(Monthly Rate, 2008 - Present)

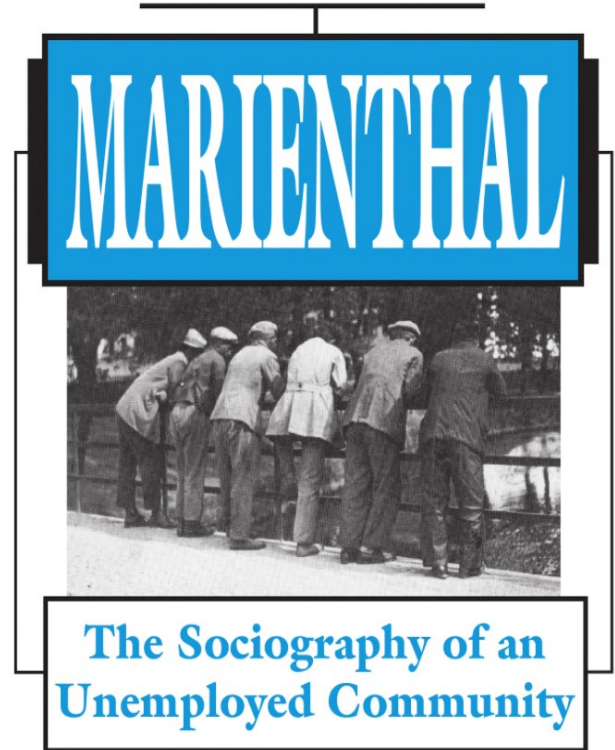
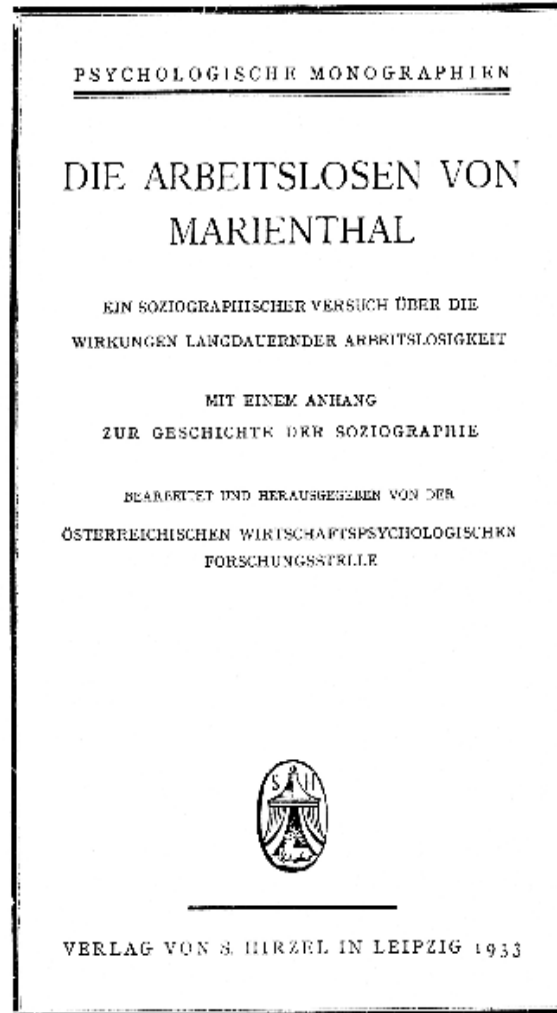
- Official CPI
- PriceStats Index



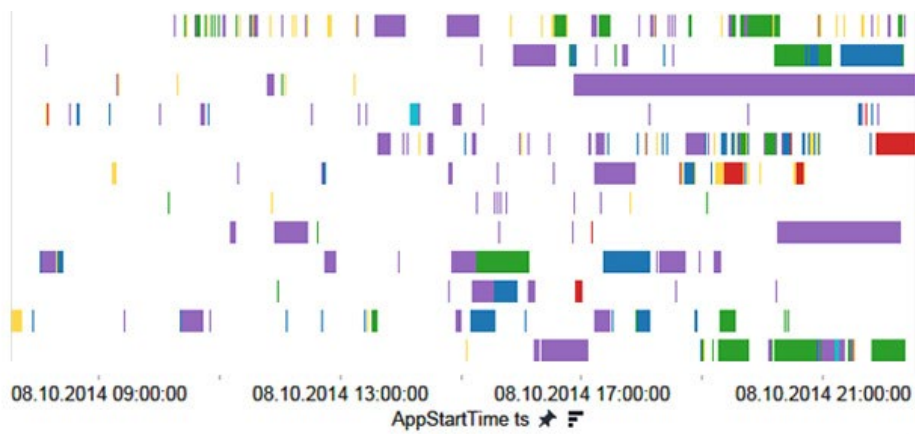
Source: State Street, PriceStats

# Observations

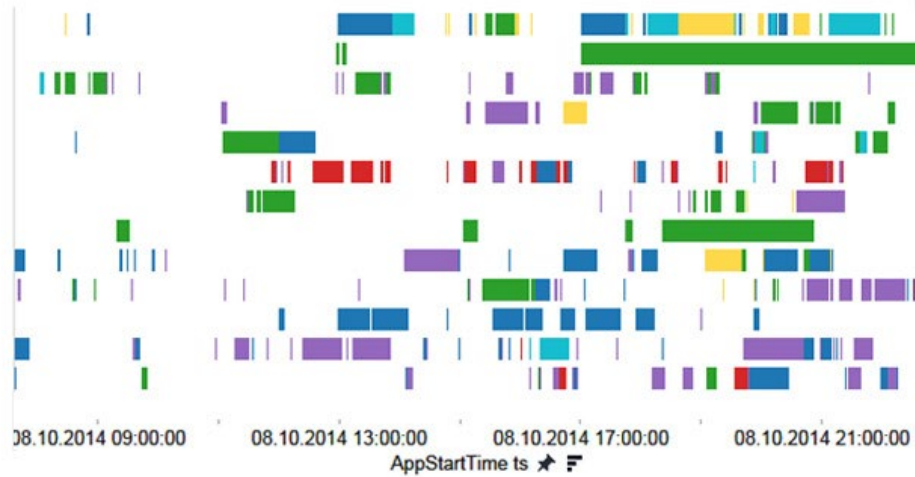
**GitHub**



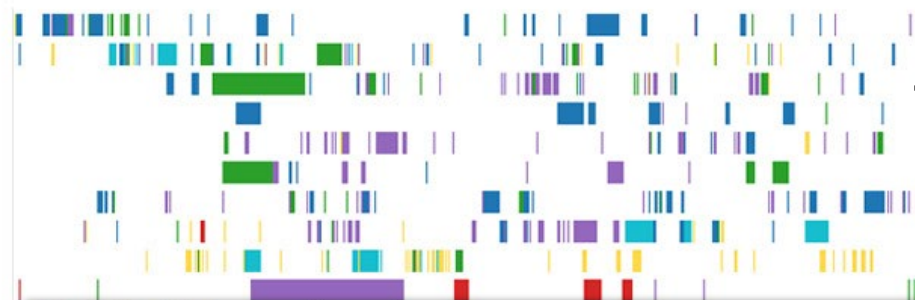
Marie Jahoda, Paul F. Lazarsfeld,  
and Hans Zeisel



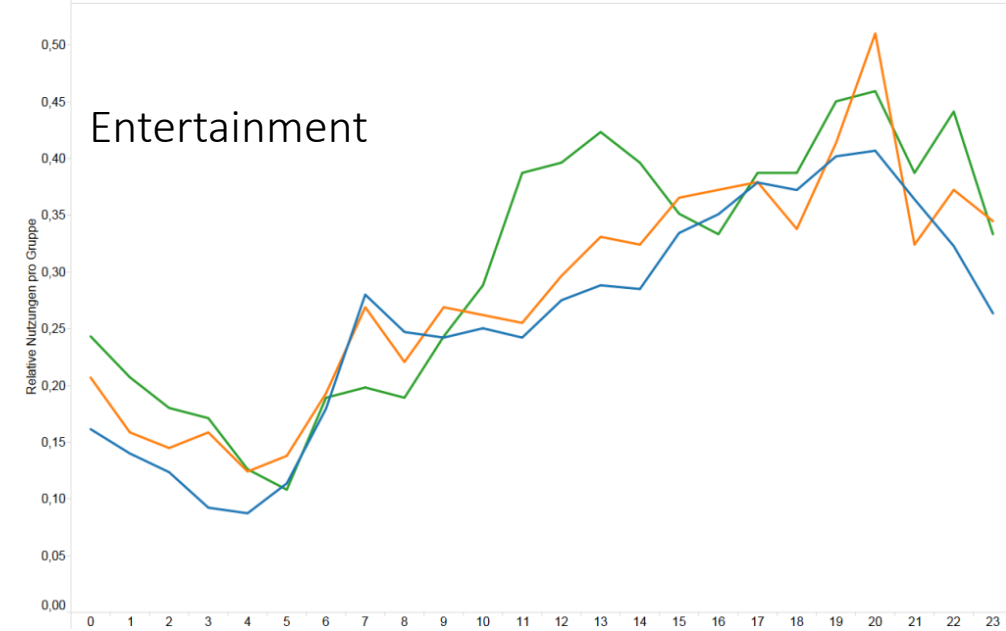
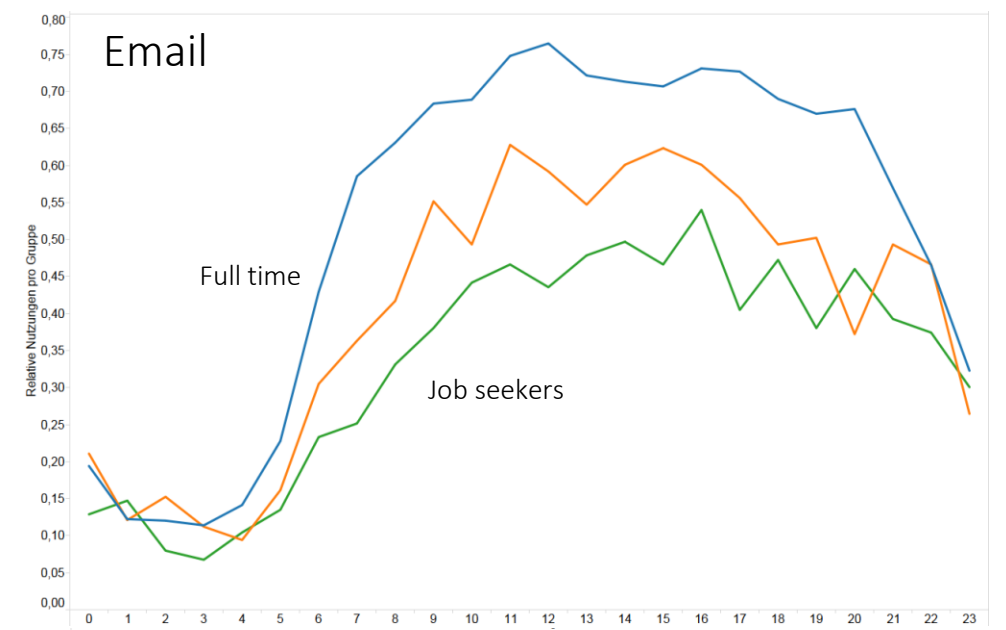
Full-time employed  
→ App use past 5pm



Part-time employed  
→ App use at noon



Job seekers  
→ Continuous app use



However ...

Work & \$

Experiment

Survey

2  
3 Steuernummer  
4  
5 Einkünfte aus nichtselbständiger Arbeit  
6 Angaben zum Arbeitslohn  
7 Steuerklasse 100  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100

Bruttoarbeitslohn	110
Lohnsteuer	140
Sonderausgaben	400





Experiment

Survey

A screenshot of a German tax form titled "Einkünfte aus nichtselbständiger Arbeit". The form includes a tax number field, a section for "Angaben zum Arbeitslohn" (Details on wages), and a table with the following data:

	Steuernummer	Steuernummer
5	Steuernummer	Steuernummer
6	Bruttobehalt	110
7	Lohnsteuer	140
8	Sonderausgaben	400



Work & \$

New Skills

Data Output/Access

Data Analysis

Data Curation/Storage

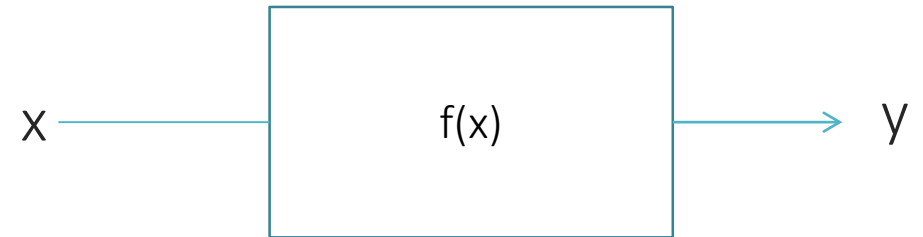
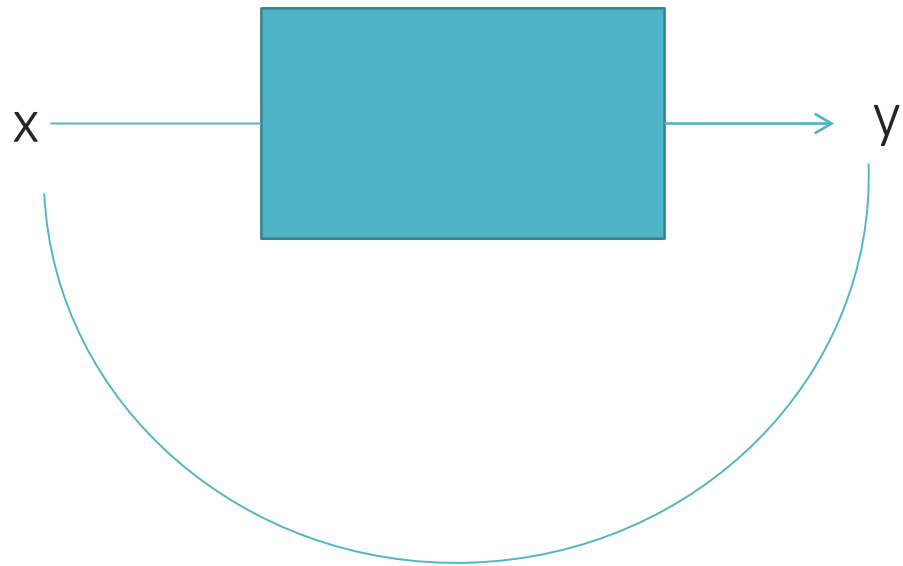


Learn how to communicate results and distribute and store your data

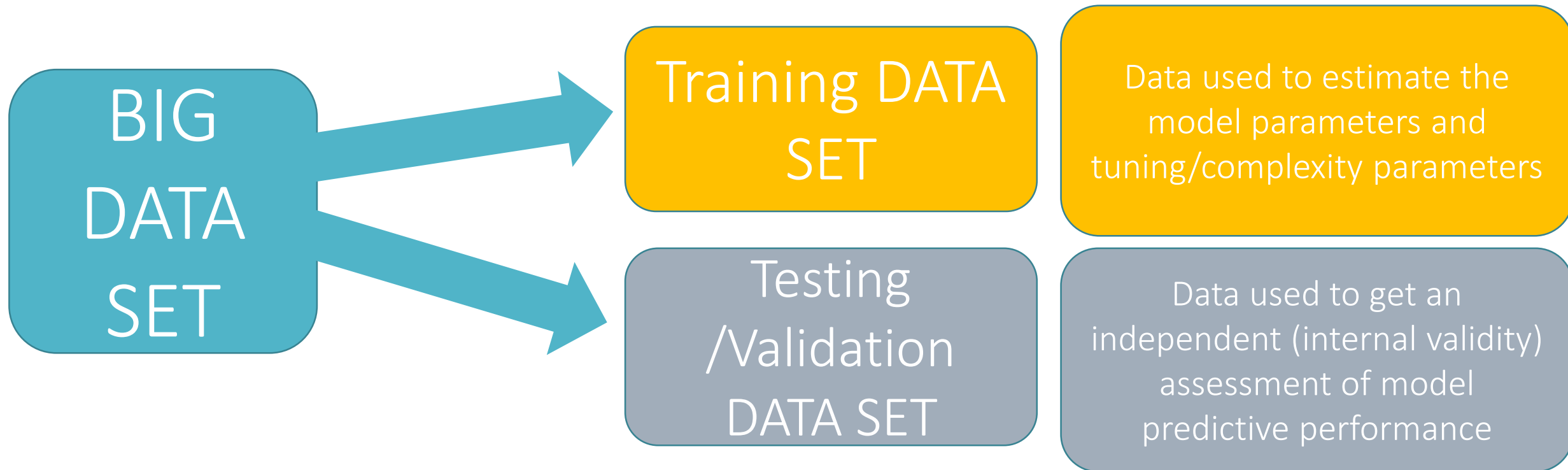
Learn a variety of analysis methods suited for different data types

Learn how to curate and manage data

# Machine Learning - AI



# Model Evaluation Strategy: Split Sample



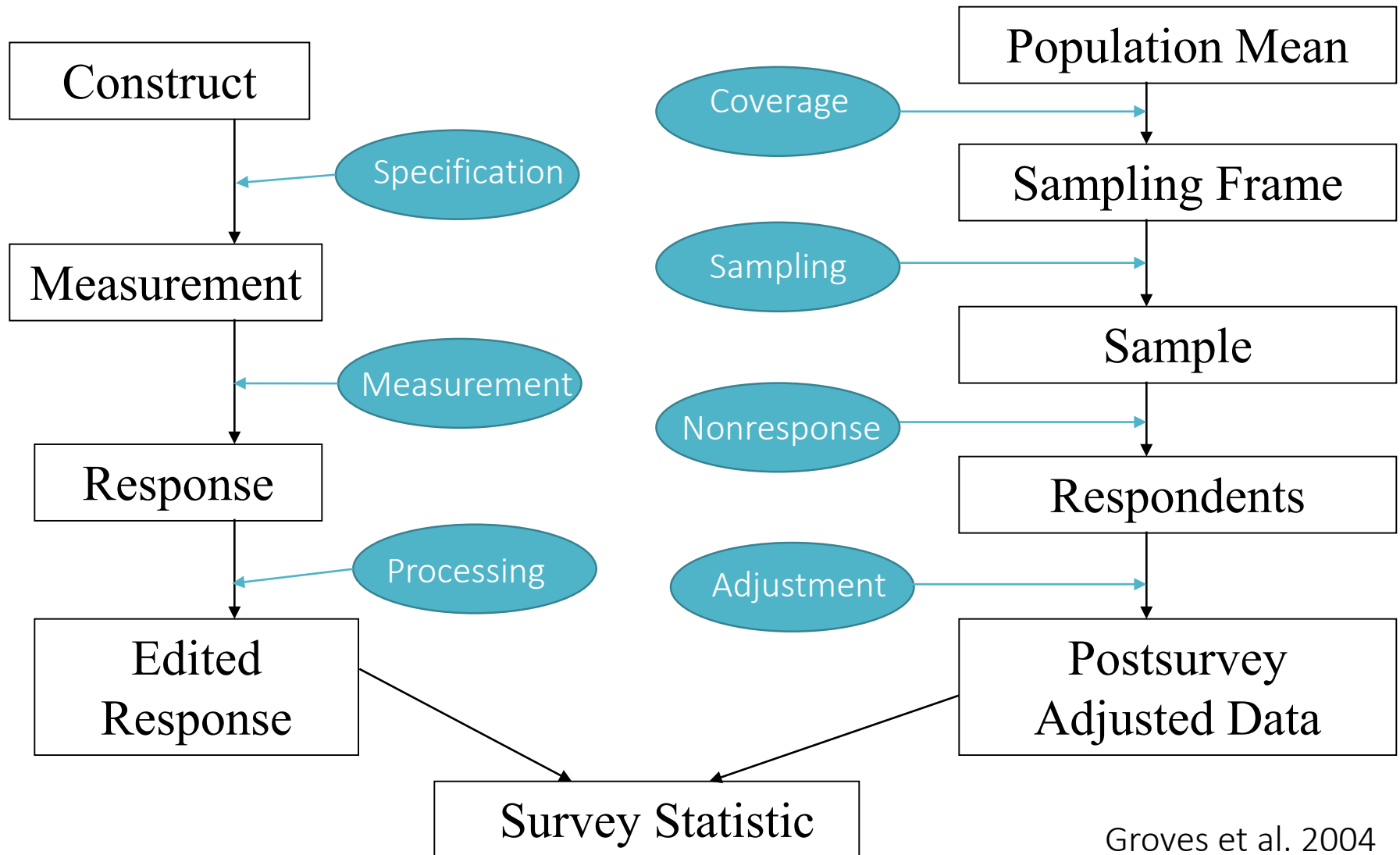
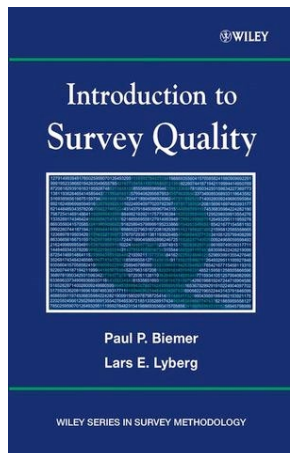
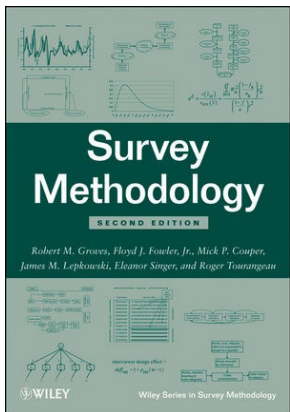


# Risk in Learning from Data

By far the worst approach is to wait for data quality problems to surface on their own.

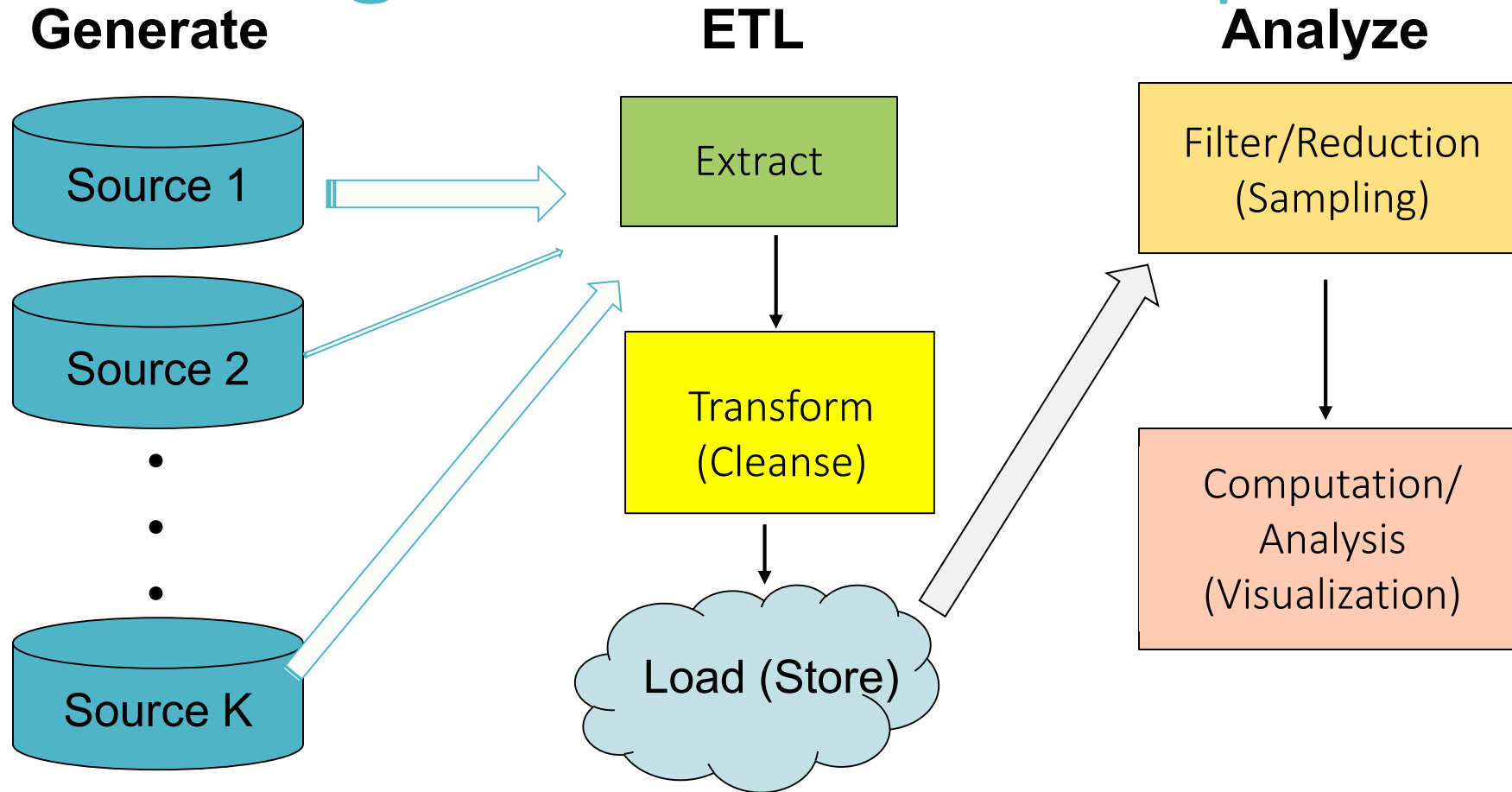
T. Herzog, F. Scheuren, W. Winkler, 2007

# Data Generating Process

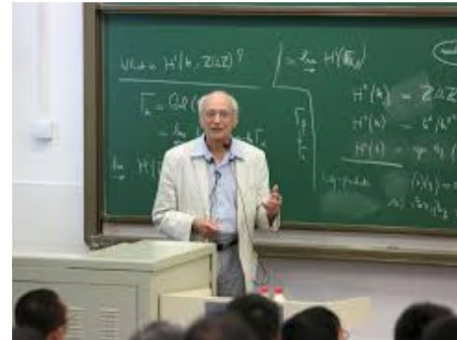
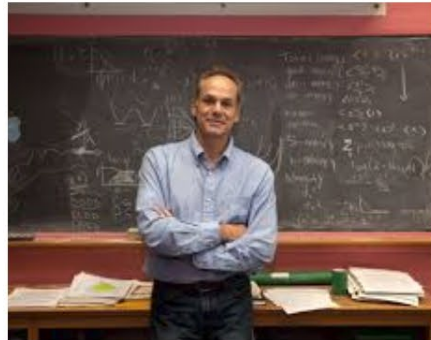
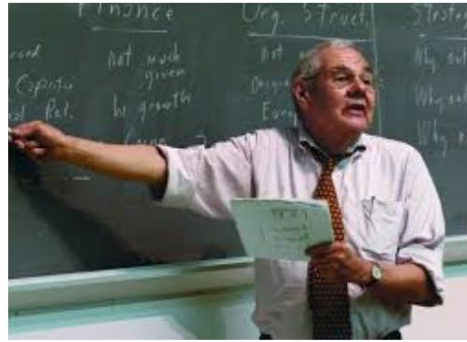


Groves et al. 2004

# Big Data Process Map



Source: Paul Biemer in Japac, Kreuter et al. 2015 – AAPOR Task Force Report



Google Image Search: June 3<sup>rd</sup> 2018

Search Term: "University Professor"





## DOMAIN EXPERT

User, analyst, or leaders with deep subject matter expertise related to the data, its appropriate use, and its limitations

## SYS ADMIN

Team member responsible for defining and maintaining a computation infrastructure that enables large scale computation

## METHODOLOGIST

Team member with experience applying formal research methods, including survey methodology and statistics

## COMPUTER SCIENTIST


Technically skilled team member with education in computer programming and data processing technology

Data Output/Access

Data Analysis

Data Curation/Storage

Data Generating Process



Learn how to communicate results and distribute and store your data

Learn a variety of analysis methods suited for different data types

Learn how to curate and manage data

Understand how to collect data yourself, and how data are generated through administrative and other processes.


Data Output/Access

Data Analysis

Data Curation/Storage

Data Generating Process

Research Question



Learn how to communicate results and distribute and store your data

Learn a variety of analysis methods suited for different data types

Learn how to curate and manage data

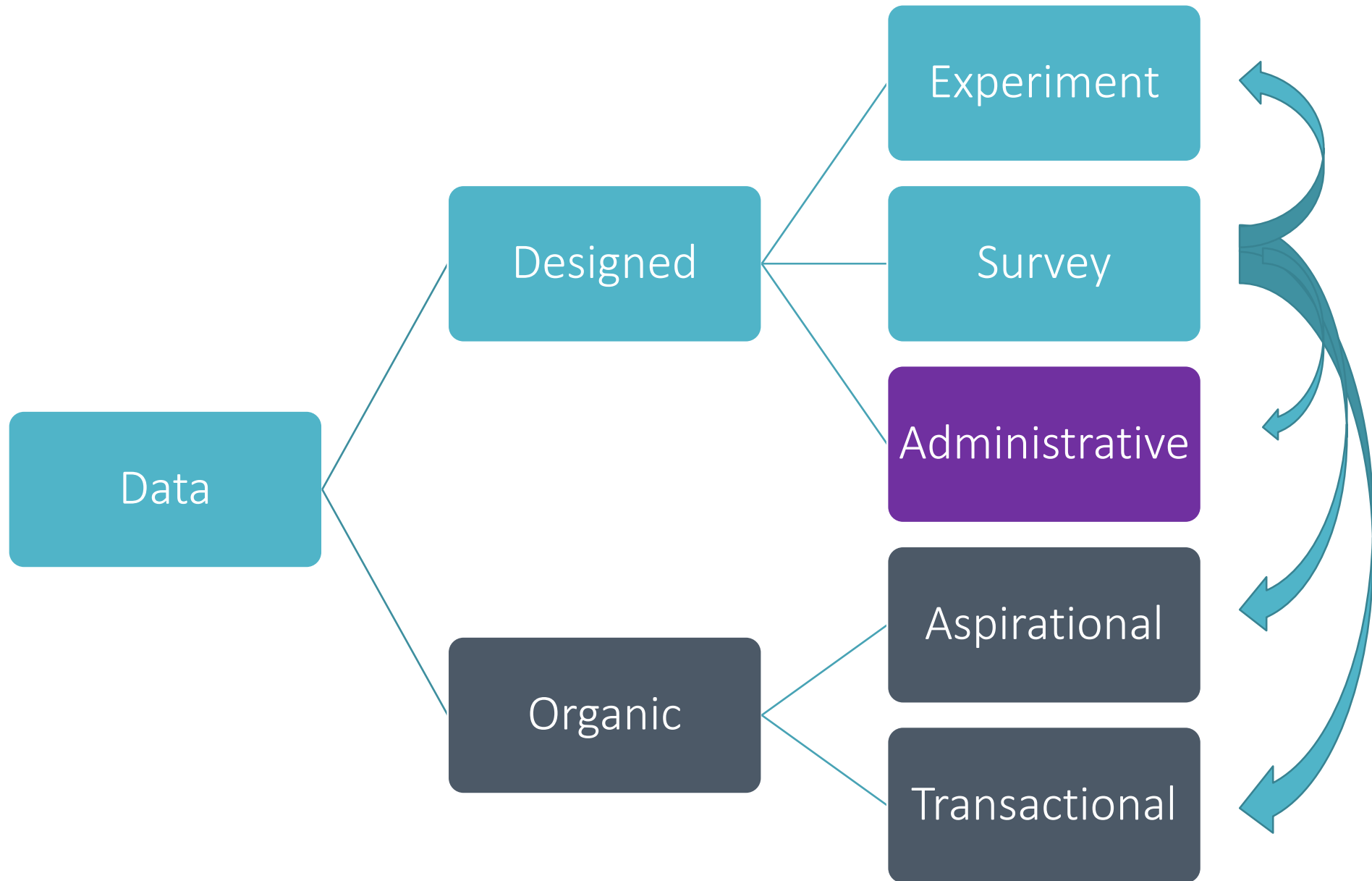
Understand how to collect data yourself, and how data are generated through administrative and other processes.

Learn how to formulate your research goal and which data are best suited to achieve this goal.

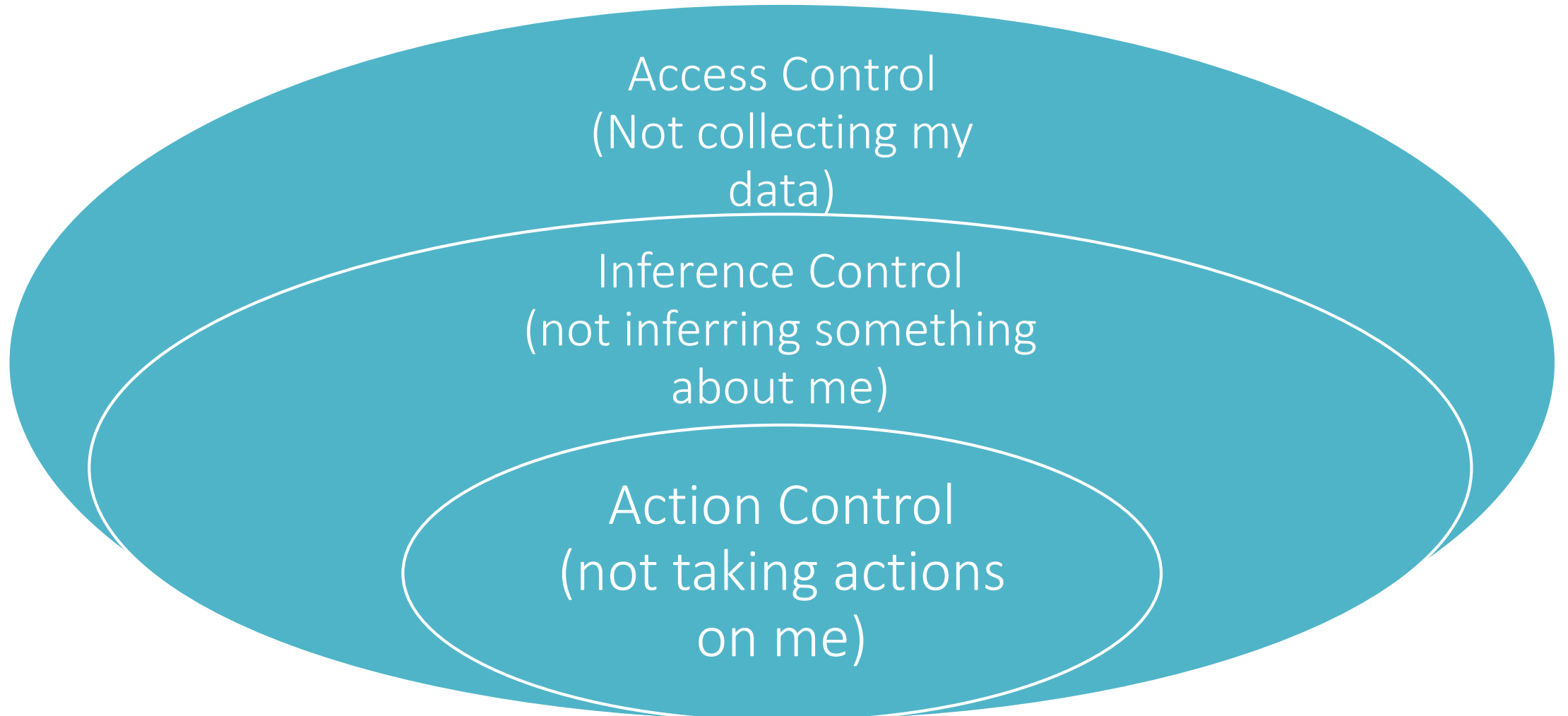
One way to think about a data analysis is to think of it as a product to be designed. [...] Producing a useful product requires careful consideration of who will be using it.

Roger Peng, 2018

# Risk in Combining Data



# Consent to give up control





### *Front*

*The data you are about to provide to us would be much more valuable if you would allow us to link them with .... Do you agree to the linkage?*

### *Back*

*The data you provided to us would be much more valuable if you would allow us to link them with .... Do you agree to the linkage?*

Phone	Front	Back	Total n
% agree	90.8	78.7	598

Web	Front	Back	Total
% agree	82.6	62.4	520

One Option

# Introduction to International Program in Survey and Data Science



SPONSORED BY THE



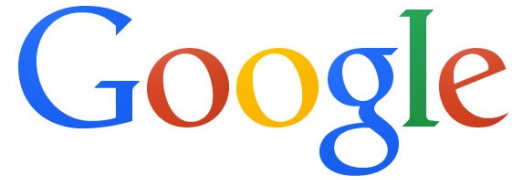
Federal Ministry  
of Education  
and Research



# International Faculty from Partner Universities



# International Faculty from the Industry



Institute for Employment  
Research

The Research Institute of the  
Federal Employment Agency





# Selection of Courses

- Analysis of Complex Survey Data, 2 cr./4 ECTS
- Big Data and Machine Learning, 1 cr./2 ECTS
- Computer-Based Content Analysis I, 1 cr./2 ECTS
- Computer-Based Content Analysis II, 1 cr./2 ECTS
- Data Collection Methods, 3 cr./6 ECTS
- Experimental Design for Surveys, 2 cr./4 ECTS
- Fundamentals of Survey and Data Science, 3 credits/6 ECTS
- Generalized Linear Models, 2 cr./4 ECTS
- Inference from Complex Surveys, 2 cr./4 ECTS
- Introduction to Data Visualization, 1 cr./2 ECTS
- Introduction to Python and SQL, 1 cr./2 ECTS
- Introduction to Real World Data Management, 2 cr./4 ECTS
- Introduction to Small Area Estimation, 2 cr./4 ECTS
- Practical Tools for Sampling & Weighting , 2 cr./4 ECTS
- Privacy Law, 1 cr./2 ECTS
- Project Consulting, 6 cr./12 ECTS
- Questionnaire Design, 2 cr./4 ECTS
- Review of Statistical Concepts (bridge course)
- Web Survey Methodology, 2 cr./4 ECTS



# Open Enrollment Courses - 2019/2020

- Applied Sampling (Sampling I), 2 credits/4 ECTS
- Data Confidentiality and Statistical Disclosure Control, 2 credits/4 ECTS
- Introduction to Record Linkage with Big Data Application, 1 credit/2 ECTS
- Item Nonresponse and Imputation, 1 credit/2 ECTS
- Measurement Error Models, 1 credit/2 ECTS
- Usability Testing for Survey Research, 1 credit/2 ECTS
- Web Scraping and API, 1 credit/2 ECTS

# Flexible & engaging online learning environment

- Online learning environment accessible from anywhere in the world (taught in English)
  - Small virtual classroom with a mix of synchronous & asynchronous learning
    - *Pre-recorded lectures split into small video units*
    - *Required readings and (bi)weekly assignments*
    - *Discussion forums*
    - *Weekly online meetings*
- 8-10h per week
- Annual on-site networking activity with fellow students from five continents
  - Wide variety of options: from individual courses or course sequences to a modular program

# Format

## Asynchronous

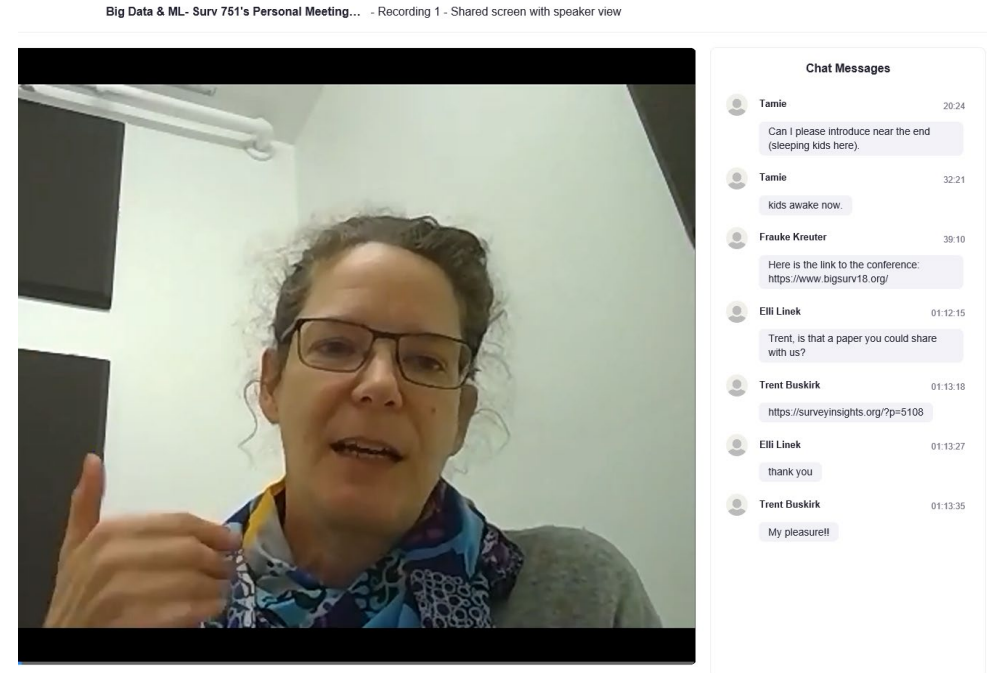


Dr. Thomas Fetzer, LL.M. (Vanderbilt)



- Pre-recorded lectures (split into small video units)
- Required readings and (bi)weekly assignments
- Discussion forums

## Synchronous



- Small virtual classrooms
- Weekly 50-minute discussions led by the instructor

# Summary

1. Data are plentiful and increasingly accessible
2. Costs of statistic production often not reduced but shifted from collection to processing
3. New data need new skills. Those can be acquired by existing staff, manageable by social scientists and content experts. Best to work in teams
4. Mindset shift to focus on quality (as it is strength of UN) is even more important
5. Ethical and cognitive psychology consideration important when (asking for) combining data, to stay within a contextual integrity of data use

THANK YOU!

[fkreuter@umd.edu](mailto:fkreuter@umd.edu)